

50277-0378

Patent

UNITED STATES PATENT APPLICATION

FOR

METHOD AND APPARATUS FOR DEBUGGING A SOFTWARE PROGRAM USING DYNAMIC
DEBUG PATCHES AND COPY ON WRITE VIEWS

INVENTORS:

VIKRAM JOSHI
ALEX TSUKERMAN
SHARI YAMAGUCHI

PREPARED BY:

HICKMAN PALERMO TRUONG & BECKER, LLP
1600 WILLOW STREET
SAN JOSE, CA 95125-5106
(408) 414-1080

EXPRESS MAIL CERTIFICATE OF MAILING

"Express Mail" mailing label number EL652871070US

Date of Deposit November 20, 2000

I hereby certify that this paper or fee is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 CFR 1.10 on the date indicated above and is addressed to the Assistant Commissioner for Patents, Washington, D.C. 20231.

Casey Moore

(Typed or printed name of person mailing paper or fee)

Casey Moore

(Signature of person mailing paper or fee)

METHOD AND APPARATUS FOR DEBUGGING A SOFTWARE PROGRAM
USING DYNAMIC DEBUG PATCHES AND COPY ON WRITE VIEWS

CROSS REFERENCE TO RELATED APPLICATIONS

This application is related to and claims domestic priority under 35 U.S.C. § 119(e) from prior U.S. Provisional Patent Application Serial Number 60/166,598 filed on November 19, 1999 entitled "Debugging Techniques And Fast SGA Dumps For Deferred
5 Analysis Of The Database", by inventors Vikram Joshi, Alex Tsukerman, and Shari Yamaguchi, the entire disclosure of which is hereby incorporated by reference as if fully set forth herein.

This application is related to U.S. Patent Application Serial Number 09/649,310 filed on August 28, 2000, entitled "Method And Apparatus For Debugging A Software
10 Program", by inventors Vikram Joshi, Alex Tsukerman, and Shari Yamaguchi, the entire disclosure of which is hereby incorporated by reference as if fully set forth herein.

This application is related to U.S. Patent Application Serial Number 09/717,162,
filed on the same day herewith entitled "Fast Database State Dumps To File For Deferred
Analysis Of A Database", by inventors Vikram Joshi, Alex Tsukerman, and Shari
15 Yamaguchi (Attorney docket number 50277-379, OID 1999-173-03), the entire disclosure of which is hereby incorporated by reference as if fully set forth herein.

This application is related to U.S. Patent Application Serial Number 09/717,161,
filed on the same day herewith entitled "A Debug And Data Collection Mechanism
Utilizing A Difference In Database State By Using Consecutive Snapshots Of A
20 Database State", by inventors Vikram Joshi, Alex Tsukerman, and Shari Yamaguchi
(Attorney docket number 50277-380, OID 1999-173-04), the entire disclosure of which is hereby incorporated by reference as if fully set forth herein.

FIELD OF THE INVENTION

The present invention generally relates to debugging software programs and, more specifically, to techniques for debugging database systems.

BACKGROUND OF THE INVENTION

5 In a database system, an area of system memory is allocated and one or more processes are started to execute one or more transactions. The database server communicates with connected user processes and performs tasks on behalf of the user. These tasks typically include the execution of transactions. The combination of the allocated system memory and the processes executing transactions is commonly termed a database "server" or "instance".

10 Like most software systems, a database server has complicated shared memory structures. A shared memory structure contains data and control information for a portion of a database system. Because of software, hardware, or firmware bugs that may exist in a complex database system, shared memory structures may become logically incorrect. When structures become logically incorrect, the database system is likely to fail. Database failure is typically discovered in the following ways: by checking consistency of structures; by verifying certain assumptions; or by running into corrupted pointers. Attempting to process corrupted pointers will lead to a "crash," after which normal database operation is no longer possible.

15 A major responsibility of the database administrator is to be prepared for the possibility of hardware, software, network, process, or system failure. When shared structures are presumed to be corrupted, the best course of action for a database administrator is to cease further processing of the database. If a failure occurs such that the operation of a database system is affected, the administrator must usually recover the database and return the database to normal operations as quickly as possible. Recovery should protect the database and associated users from unnecessary problems and avoid or reduce the possibility of having to duplicate work manually.

Recovery operations vary depending on the type of failure that occurred, the structures affected, and the type of recovery that is performed. If no files are lost or damaged, recovery may amount to no more than rebooting the database system. On the other hand, if data has been lost, recovery requires additional steps in order to put the database back into normal working order.

Once the database is recovered or rebooted, the immediate problem is quickly resolved, but because the root cause is still undetermined and therefore unresolved, the error condition may resurface, potentially causing several additional outages. Therefore, it is still important to diagnose the state of the structures and data surrounding the database failure. Such a diagnosis may provide valuable information that can reduce the chance of failure in the future. As a practical matter, diagnosing the failure may lead to determining which vendor's hardware or software is responsible for the database failure. Such information is valuable for a vendor's peace of mind, if nothing else. Thus, competing with the goal of recovering the database as quickly as possible, is the goal of determining why the database system failed in the first place.

Unfortunately, even with traditional techniques of diagnosing a database failure, the system administrator is usually unable to obtain a sufficient amount of clues to determine why the failure happened. A deliberate and thorough diagnosis of the failure may require an unacceptable amount of database downtime. For example, any amount of downtime over 30 minutes may be extremely costly for a database that is associated with a highly active web site. Too much downtime may have unduly expensive business ramifications, such as lost revenue and damage to the reputation of the web site owner.

Another problem with traditional debugging techniques is that they can be intrusive. For example, a database system that supports the Structured Query Language (SQL) may be debugged by compiling SQL statements and running them against the database. The act of compiling and executing the SQL statements changes the state of

database server's end-memory state, which is the state after the database has been shutdown. Because the end-memory state is being analyzed separately from the database, the programmer performing the debugging does not have access to the real database and some of the database's persistent structures. Some of these persistent structures could be on disk or, in a multiple node system, on other nodes. For example, in a parallel server configuration, the persistent structures needed for debugging could reside on other servers. Thus, the technique of separately debugging portions of the database prevents the programmer from using the data that can only be obtained from the database itself.

Further, where debug operations are performed on the database while the database is down, multiple programmers cannot each privately diagnose the failure. Rather, the key data structures are typically diagnosed by having one programmer in front of a console inputting debug commands, while other programmers gather around issuing advice. Multiple programmers individually debugging the database is unadvisable using existing debugging techniques because the act of inputting debug commands is intrusive, as mentioned above. Each programmer's work would interfere with the concurrent debugging progress of fellow programmers.

For the foregoing reasons, what is needed is a method of debugging a software program, such as a database system, that is non-intrusive, yet allows for a comprehensive assessment of a failure.

SUMMARY OF THE INVENTION

Techniques are provided for allowing multiple persons to concurrently test software patches on a software program or debug a problem of the software program. Each person preferably has their own private view, which consists of (1) copied portions of the software program that reflect modifications made by that person, and (2) the portions of the preserved software program that the person has not modified. Providing a private view to each person allows each person to test and debug privately, independently, and concurrently with others. Each private view may be extensively explored and modified without affecting the memory state of the software program that existed at the time the software program was shutdown.

Accordingly, at any time, each private view may be refreshed to the state of the software program that existed at the time of shutdown. Faster diagnosis of the problem may therefore be accomplished because a debugger does not have to peek cautiously and slowly into the inner-workings of the software program. Similarly, the testing of various potential solutions to bugs in the program may be accomplished efficiently without affecting the memory state of the software program that existed at the time the software program was shutdown. Thus, where downtime of a software program must be kept to a minimum, the present techniques allow for performing quick and comprehensive diagnostics and testing of potential solutions to problems in the software program.

BRIEF DESCRIPTION OF THE DRAWINGS

The present invention is illustrated by way of example, and not by way of limitation, in the figures of the accompanying drawings and in which like reference numerals refer to similar elements and in which:

5 FIG. 1 illustrates how data is copied to preserve the memory state of a software program before being modified in response to a technique for debugging and testing of potential solutions for a software program;

FIG. 2 is a flowchart of a technique that allows for non-intrusive debugging and testing of potential solutions for a software program; and

10 FIG. 3 is a block diagram that illustrates a computer system upon which an embodiment of the invention may be implemented.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

Techniques for non-intrusive testing of potential solutions for and debugging of a software program are described. In the following description, for the purposes of explanation, numerous specific details are set forth in order to provide a thorough understanding of the present invention. It will be apparent, however, to one skilled in the art that the present invention may be practiced without these specific details. In other instances, well-known structures and devices are shown in block diagram form in order to avoid unnecessarily obscuring the present invention.

TESTING AND DIAGNOSTIC TECHNIQUE

A database administrator may cause a database system to cease execution for a number of reasons, which are discussed above. Diagnosing a database will typically lead to modifying data in the database while it is down. As explained above, it is desirable to preserve the memory state that existed at the time of failure or at the time the database was shutdown.

FIG. 1 illustrates how, using the techniques described herein, data is copied to preserve the memory state of a software program before being modified. For the purpose of explanation, it will be assumed that the software program is a database server. However, the present techniques are not limited to any particular type of program. A database administrator, for example, "freezes" a portion of a database to preserve the memory state of the database system. Preserving the database may include suspending a failed process within the database system. Various techniques may be used to freeze the state of a database server. One such technique is described in U.S. Patent Application No. 09/223,660 entitled "METHOD AND SYSTEM FOR DIAGNOSTIC PRESERVATION OF THE STATE OF A COMPUTER SYSTEM" filed by Wei Hu and Juan Loaiza on December 30, 1998, the contents of which is incorporated herein by reference.

The act of preserving the database may be initiated by giving the software

program an explicit "freeze" user command. Alternatively, the act of preserving the database may be initiated in response to an automatic trigger that fires when an error event is detected. The techniques described hereafter indicate how debugging and testing may be performed, even by multiple users concurrently, without changing any data in the preserved portion 102. Any operation that could cause any change within the preserved portion is disabled with respect to the preserved portion 102.

In one embodiment, the software program is a database server that is composed of a memory portion referred to as a "database instance", and a set of data on disk referred to as "datafiles". In addition, one database server may be shared in a hardware cluster with additional database instances residing on different nodes, but still sharing access to the same set of datafiles. Preferably, the database administrator will be able to issue a command to preserve only the database instance that has failed, thereby detaching the preserved database instance from the cluster membership. In such a situation, it is important to exercise care while using the preserved database instance in this state, and to not affect the integrity of the datafiles or other database instances. Preferably, a database administrator only uses such a detached preserved database instance for debugging and testing operations. Detaching a database instance from a cluster involves isolating the instance from the rest of the cluster using software and/or hardware means.

SECONDARY SOFTWARE PROGRAMS

While the software program is preserved, debugging and or testing operations may begin. Debugging or testing may involve a second software program. The second software program may be a software patch that fixes a bug in the software program. Techniques are provided to allow for testing of the software patch without compromising the memory state of the original software program. Alternatively, the second software program may be a diagnostic tool for debugging the software program. Any such second

software program is herein generally referred to as a "secondary" software program.

FIG. 1 shows the preserved portion 102 having segments of data. These segments of data are typically pages of memory. Assume that a user has written a secondary software program that is either a diagnostic tool or a potential patch the software program. The user may then compile and dynamically link the secondary software program to the original software program. Many operating systems provide for the compilation and dynamically linkage of separate software programs.

Assume that when the user executes the secondary software program, segments of the preserved portion 102 are accessed. For example, execution of the secondary results in a read operation that accesses data 104. The execution of the second software program may also call for a write operation to be performed on some data within the preserved portion 102. The data within the reserved portion that is targeted by the write operation is referred to herein as targeted data 106a. In response to an attempt to perform a write operation on data within the targeted data 106a, a copy is made of the targeted data 106a. The actual modification that would have been made to the targeted data 106a is instead made to the copy, creating a modified copy 106b. In one embodiment, the modified copy 106b is a copy-on-write page of memory.

In subsequent operations, relative to the execution of the software program, the modified copy 106b takes the place of the targeted data 106a. Thus, if the execution of the secondary software program involves a subsequent read operation of targeted data 106a, the read operation would be performed on the modified copy 106b. Similarly, if the execution of the secondary software program would involve further modification to the targeted data 106a, the modification would once again be performed on the modified copy 106b.

For simplification purposes, FIG. 1 shows a scenario in which the target portion 206a includes only one segment. Therefore, a modified copy 106b has been made of

only one target portion 106a. However, the execution of the secondary software program may actually involve modifications to many areas of a software program, and therefore cause the generation of modified copies of a multitude of segments.

In one embodiment, the memory segments of the preserved portion 102 are pages in memory. Preferably, a page map is used to keep track of all the pages of the modified copies for each user. For example, the modified copy 106b may be a copy-on-write page, the address of which is kept in a page map. Using this copying technique, the original preserved portion 102 of the software program is unaltered. The page mapping software and/or hardware ensures that, for the user that caused creation of the modified page, the modified page is mapped at the same virtual address as the original page, thus preserving the integrity of data structure references, indices, and pointers. The modified page is a modified copy of the original page, and the user that caused the creation of the modified page has a private view of the modified page. Such modified pages are herein referred to as "view private" modified pages.

The modified copies for a user may be discarded at anytime. Thus, a fresh testing and debug session may be initiated at any time using the preserved portion 102 and the aforementioned copying technique.

MAINTAINING PER-USER MODIFICATION DATA

According to one embodiment, a separate set of modified copies are maintained for each user, based on the modifications made by that user. Specifically, each user sees (1) the modified copies that have been generated in response to the execution of secondary software programs executed by that user, and (2) the preserved portions of the software program that have not been modified in response to the execution of secondary software programs initiated by that user. The modified copies may be managed in a view private fashion using any one of a number of techniques, including page mapping software and hardware techniques.

Consequently, multiple programmers may debug and test potential solutions to problems in the software program concurrently, independently, and privately. The debug and testing progress of one programmer will not affect the debug and testing progress of another programmer. Accordingly, any number of debug and testing sessions may be generated and later destroyed. Such multiple-session debugging and testing should lead to a relatively quick and comprehensive assessment of the failure and testing of various potential solutions to the failure. For example, in FIG. 1, in response to compiling, dynamically linking and executing software program 1, which is a secondary software program, a user 1 modifies 106a and 106b is created. In response to compiling, dynamically linking and executing software program 2, which is another secondary software program that is distinct from software program 1, a user 2 modifies 106a and 106c is created. User 1 is unable to see the changes in 106c, and user 2 is unable to see the changes in 106b.

SHARING SECONDARY PROGRAMS

In addition, if user 1 wishes to share software program 1, user 1 may provide access to software program 1 by publishing a symbolic name associated with software program 1 in slot 108 of FIG. 1. Multiple users may concurrently, independently, and privately execute software program 1 by using the corresponding symbolic name that is published in the preserved portion of the software program. For example, in response to user 3 concurrently, independently, and privately executing software program 1, 106d in FIG. 1 is created. Because modified copies may be managed in a view private fashion, user 3 is unable to see the changes in either 106b or 106c.

MULTIPLE SECONDARY PROGRAMS IN A SINGLE SESSION

One debug and testing session can be used to execute various secondary software programs that are separate from the original software program, which is being debugged.

example, in a parallel server configuration, the persistent structures needed for debugging could reside on other servers. The execution of a secondary software program may call for accessing data, such as a persistent structure, that is outside the current database instance being debugged. For such data that resides outside the database instance being
5 debugged, it is preferable to mount the data in a read-only mode into the database prior to performing debug and testing operations. Such a mounting step facilitates reads from outside data, such as persistent database tables. Debug and testing operations cannot write to the read-only data, and therefore will not make changes to the original persistent structures. Further, when any operation attempts to write to the read-only data, the debug
10 and testing system preferably produces a logical error message, which the debug and testing system makes known to the user.

After copying outside data into the database in a read-only mode, outside data may be treated as part of the preserved portion 102. That is, a user is allowed to perform operations that modify the data, but those operations cause the creation of separate
15 modified copies, and leave the original data intact. Thus, an embodiment of the present invention is applicable to debugging a database and testing of potential solutions to problems in the database where execution of secondary software programs call for accessing data that is outside the current database instance.

THE DEBUGGING OPERATION

20 FIG. 2 is a flowchart of a debug and testing technique that allows for non-intrusive debugging of a software program as well as testing potential solutions to problems in the software program. At block 202, a memory state of a preserved portion of the software program is preserved. As mentioned in the discussion with reference to FIG. 1, preserving the memory state may include suspending an application that has
25 failed. At block 204, a secondary software program is compiled and dynamically linked to the original software program. At block 206, the secondary software program is

executed causing the creation of a copy of targeted data 106a, if the execution of the secondary software program would normally cause modification to targeted data 106a in the portion of the software program that is being preserved.

HARDWARE OVERVIEW

5 FIG. 3 is a block diagram that illustrates a computer system 300 upon which an embodiment of the invention may be implemented. Computer system 300 includes a bus 302 or other communication mechanism for communicating information, and a processor 304 coupled with bus 302 for processing information. Computer system 300 also includes a main memory 306, such as a random access memory (RAM) or other dynamic storage device,
10 coupled to bus 302 for storing information and instructions to be executed by processor 304. Main memory 306 also may be used for storing temporary variables or other intermediate information during execution of instructions to be executed by processor 304. Computer system 300 further includes a read only memory (ROM) 308 or other static storage device coupled to bus 302 for storing static information and instructions for processor 304. A
15 storage device 310, such as a magnetic disk or optical disk, is provided and coupled to bus 302 for storing information and instructions.

Computer system 300 may be coupled via bus 302 to a display 312, such as a cathode ray tube (CRT), for displaying information to a computer user. An input device 314, including alphanumeric and other keys, is coupled to bus 302 for communicating information
20 and command selections to processor 304. Another type of user input device is cursor control 316, such as a mouse, a trackball, or cursor direction keys for communicating direction information and command selections to processor 304 and for controlling cursor movement on display 312. This input device typically has two degrees of freedom in two

axes, a first axis (e.g., x) and a second axis (e.g., y), that allows the device to specify positions in a plane.

The invention is related to the use of computer system 300 for implementing the techniques described herein. According to one embodiment of the invention, those

5 techniques are implemented by computer system 300 in response to processor 304 executing one or more sequences of one or more instructions contained in main memory 306. Such instructions may be read into main memory 306 from another computer-readable medium, such as storage device 310. Execution of the sequences of instructions contained in main memory 306 causes processor 304 to perform the process steps described herein. In

10 alternative embodiments, hard-wired circuitry may be used in place of or in combination with software instructions to implement the invention. Thus, embodiments of the invention are not limited to any specific combination of hardware circuitry and software.

The term "computer-readable medium" as used herein refers to any medium that participates in providing instructions to processor 304 for execution. Such a medium may
15 take many forms, including but not limited to, non-volatile media, volatile media, and transmission media. Non-volatile media includes, for example, optical or magnetic disks, such as storage device 310. Volatile media includes dynamic memory, such as main memory 306. Transmission media includes coaxial cables, copper wire and fiber optics, including the wires that comprise bus 302. Transmission media can also take the form of acoustic or light
20 waves, such as those generated during radio-wave and infra-red data communications.

Common forms of computer-readable media include, for example, a floppy disk, a flexible disk, hard disk, magnetic tape, or any other magnetic medium, a CD-ROM, any other optical medium, punchcards, papertape, any other physical medium with patterns of holes, a

RAM, a PROM, and EPROM, a FLASH-EPROM, any other memory chip or cartridge, a carrier wave as described hereinafter, or any other medium from which a computer can read.

Various forms of computer readable media may be involved in carrying one or more sequences of one or more instructions to processor 304 for execution. For example, the instructions may initially be carried on a magnetic disk of a remote computer. The remote computer can load the instructions into its dynamic memory and send the instructions over a telephone line using a modem. A modem local to computer system 300 can receive the data on the telephone line and use an infra-red transmitter to convert the data to an infra-red signal. An infra-red detector can receive the data carried in the infra-red signal and appropriate circuitry can place the data on bus 302. Bus 302 carries the data to main memory 306, from which processor 304 retrieves and executes the instructions. The instructions received by main memory 306 may optionally be stored on storage device 310 either before or after execution by processor 304.

Computer system 300 also includes a communication interface 318 coupled to bus 302. Communication interface 318 provides a two-way data communication coupling to a network link 320 that is connected to a local network 322. For example, communication interface 318 may be an integrated services digital network (ISDN) card or a modem to provide a data communication connection to a corresponding type of telephone line. As another example, communication interface 318 may be a local area network (LAN) card to provide a data communication connection to a compatible LAN. Wireless links may also be implemented. In any such implementation, communication interface 318 sends and receives electrical, electromagnetic or optical signals that carry digital data streams representing various types of information.

000277-137-13000

Network link 320 typically provides data communication through one or more networks to other data devices. For example, network link 320 may provide a connection through local network 322 to a host computer 324 or to data equipment operated by an Internet Service Provider (ISP) 326. ISP 326 in turn provides data communication services through the world wide packet data communication network now commonly referred to as the "Internet" 328. Local network 322 and Internet 328 both use electrical, electromagnetic or optical signals that carry digital data streams. The signals through the various networks and the signals on network link 320 and through communication interface 318, which carry the digital data to and from computer system 300, are exemplary forms of carrier waves transporting the information.

Computer system 300 can send messages and receive data, including program code, through the network(s), network link 320 and communication interface 318. In the Internet example, a server 330 might transmit a requested code for an application program through Internet 328, ISP 326, local network 322 and communication interface 318. In accordance with the invention, one such downloaded application implements the techniques described herein.

The received code may be executed by processor 304 as it is received, and/or stored in storage device 310, or other non-volatile storage for later execution. In this manner, computer system 300 may obtain application code in the form of a carrier wave.

CONCLUSION

Techniques are described above for debugging a software program and for testing potential solutions to problems in the software program. In a preferred embodiment, the software program is preserved before debug and testing operations are performed on the software program. A secondary software program, which may either be a diagnostic tool

